

Towards a Mathematical Understanding of Supervised Learning: what we know and what we don't

Weinan E

Princeton University

Joint work with **Chao Ma, Stephan Wojtowytsch, Lei Wu**

Slides can be found in: www.math.princeton.edu/~weinan

Neural network-based machine learning is both very **powerful** and very **fragile**.

- What are the reasons behind?
- How can we do better (more robust formulation)?

Basic problem of supervised learning (regression)

Given $S = \{(\mathbf{x}_j, y_j = f^*(\mathbf{x}_j)), j \in [n]\}$, learn (i.e. approximate) f^* .

- assume $\mathbf{x}_j \in X = [0, 1]^d$, $\mu =$ the distribution of $\{\mathbf{x}_j\}$

Standard procedure:

- 1 choose a hypothesis space (set of trial functions) \mathcal{H}_m

- neural network models

- 2 choose a loss function (to fit the data)

- “empirical risk”

$$\hat{\mathcal{R}}_n(f) = \frac{1}{n} \sum_j (f(\mathbf{x}_j) - y_j)^2 = \frac{1}{n} \sum_j (f(\mathbf{x}_j) - f^*(\mathbf{x}_j))^2$$

- 3 choose an optimization algorithm and the hyper-parameters

- gradient descent (GD), stochastic gradient descent (SGD), ADAM, RMSprop, ...

Objective: Minimize the “population risk” (the “generalization error”)

$$\mathcal{R}(f) = \mathbb{E}_{\mathbf{x} \sim \mu} (f(\mathbf{x}) - f^*(\mathbf{x}))^2 = \int_{\mathbb{R}^d} (f(\mathbf{x}) - f^*(\mathbf{x}))^2 d\mu$$

Three important aspects to study

- hypothesis space: what kind of functions can be approximated efficiently, generalization gap (= difference between training and testing errors)
- loss function (the variational problem): landscape?
- training: can we optimize? does the solution generalize?

Important parameters:

- m =: number of free parameters
- n =: size of training dataset
- t =: training steps
- d =: dimensionality

typically interested in the case: $m, n, t \rightarrow \infty, d \gg 1$.

Main goal: Error estimates

$$\mathcal{R}(f_{m,n,t}) \lesssim m^{-\alpha} + n^{-\beta} + t^{-\gamma}$$

Want: free of **curse of dimensionality (CoD)**: α, β, γ are independent of d .

Given a class of hypothesis space (e.g. two-layer neural networks):

1. What class of functions can be approximated by that model without CoD?

Looking for some function spaces. These are the analog of Besov spaces.

2. Generalization gap for this function class?

Bounds on Rademacher complexity.

Function spaces for ML models

- Random feature model: RKHS (reproducing kernel Hilbert space,)
- Two-layer neural networks: Barron space (Bach (2017), E, Ma and Wu (2018, 2019))
- ResNets: Flow-induced space (E, Ma and Wu (2019))
- Multi-layer neural networks: Multi-layer spaces (E and Wojtowytsch (2020))

Looking for the right function spaces

- Given a type of hypothesis space \mathcal{H}_m , identify the natural function space associated with them (in particular, identify a norm Identify $\|f^*\|_*$)

- Direct approximation theorem (with Monte Carlo rate):

$$\inf_{f \in \mathcal{H}_m} \mathcal{R}(f) = \inf_{f \in \mathcal{H}_m} \|f - f^*\|_{L^2(d\mu)}^2 \lesssim \frac{\|f^*\|_*^2}{m}$$

- Inverse approximation theorem: If a function f^* can be approximated efficiently by the functions in \mathcal{H}_m , as $m \rightarrow \infty$ with some uniform bounds, then $\|f^*\|_*$ is finite.
- Study the generalization gap for this function space. One way to do this is to study the Rademacher complexity of the set $\mathcal{H}_Q = \{f, \|f\|_* \leq Q\}$.
 - Ideally, we would like to have (\hat{f} = output of ML model):

$$\text{Rad}_S(\mathcal{H}_Q) \lesssim \frac{Q}{\sqrt{n}}$$

Combined: Up to log terms, we have

$$\mathcal{R}(\hat{f}) \lesssim \frac{\|f^*\|_*^2}{m} + \frac{\|f^*\|_*}{\sqrt{n}}$$

Two-layer neural network model: Barron spaces

E, Ma and Wu (2018, 2019), Bach (2017) ⁱ

$$\mathcal{H}_m = \{f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})\}, \theta = \{(a_j, \mathbf{w}_j), j \in [m]\}$$

Consider the function $f : X = [0, 1]^d \mapsto \mathbb{R}$ of the following form

$$f(\mathbf{x}) = \int_{\Omega} a \sigma(\mathbf{w}^T \mathbf{x}) \rho(da, d\mathbf{w}) = \mathbb{E}_{(a, \mathbf{w}) \sim \rho} [a \sigma(\mathbf{w}^T \mathbf{x})], \quad \mathbf{x} \in X$$

$\Omega = \mathbb{R}^1 \times \mathbb{R}^{d+1}$, ρ is a probability distribution on Ω .

$$\|f\|_{\mathcal{B}} = \inf_{\rho \in P_f} \left(\mathbb{E}_{\rho} [a^2 \|\mathbf{w}\|_1^2] \right)^{1/2}$$

where $P_f := \{\rho : f(\mathbf{x}) = \mathbb{E}_{\rho} [a \sigma(\mathbf{w}^T \mathbf{x})]\}$.

$$\mathcal{B} = \{f \in C^0 : \|f\|_{\mathcal{B}} < \infty\}$$

ⁱRelated work in Barron (1993), Klusowski and Barron (2016), E and Wojtowytsch (2020)

Theorem (Direct Approximation Theorem)

$$\|f - f_m\|_{L^2(X)} \lesssim \frac{\|f\|_{\mathcal{B}}}{\sqrt{m}}$$

Theorem (Inverse Approximation Theorem)

Let

$$\mathcal{N}_C \stackrel{\text{def}}{=} \left\{ \frac{1}{m} \sum_{k=1}^m a_k \sigma(\mathbf{w}_k^T \mathbf{x}) : \frac{1}{m} \sum_{k=1}^m |a_k|^2 \|\mathbf{w}_k\|_1^2 \leq C^2, m \in \mathbb{N}^+ \right\}.$$

Let f^* be a continuous function. Assume there exists a constant C and a sequence of functions $f_m \in \mathcal{N}_C$ such that

$$f_m(\mathbf{x}) \rightarrow f^*(\mathbf{x})$$

for all $\mathbf{x} \in X$, then there exists a probability distribution ρ^* on Ω , such that

$$f^*(\mathbf{x}) = \int a \sigma(\mathbf{w}^T \mathbf{x}) \rho^*(da, d\mathbf{w}),$$

for all $\mathbf{x} \in X$ and $\|f^*\|_{\mathcal{B}} \leq C$.

Theorem (Bach, 2017)

Let $\mathcal{F}_Q = \{f \in \mathcal{B}, \|f\|_{\mathcal{B}} \leq Q\}$. Then we have

$$\text{Rad}_S(\mathcal{F}_Q) \leq 2Q \sqrt{\frac{2 \ln(2d)}{n}}$$

A priori estimates for regularized model

$$\mathcal{L}_n(\theta) = \hat{\mathcal{R}}_n(\theta) + \lambda \sqrt{\frac{\log(2d)}{n}} \|\theta\|_{\mathcal{P}}, \quad \hat{\theta}_n = \operatorname{argmin} \mathcal{L}_n(\theta)$$

where the path norm is defined by:

$$\|\theta\|_{\mathcal{P}} = \left(\frac{1}{m} \sum_{k=1}^m |a_k|^2 \|\mathbf{w}_k\|_1^2 \right)^{1/2}$$

Theorem (E, Ma, Wu, 2018)

Assume $f^ : X \mapsto [0, 1] \in \mathcal{B}$. There exist constants C_0 , such that for any $\delta > 0$, if $\lambda \geq C_0$, then with probability at least $1 - \delta$ over the choice of training set, we have*

$$\mathcal{R}(\hat{\theta}_n) \lesssim \frac{\|f^*\|_{\mathcal{B}}^2}{m} + \lambda \|f^*\|_{\mathcal{B}} \sqrt{\frac{\log(2d)}{n}} + \sqrt{\frac{\log(1/\delta) + \log(n)}{n}}.$$

Other models

1. Random feature model

$\{\phi(\cdot; \mathbf{w})\}$: collection of random features, e.g. $\phi(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$.

π : prob distribution of the random variable \mathbf{w} .

Hypothesis space: Given any realization $\{\mathbf{w}_j\}_{j=1}^m$, i.i.d. with distribution π

$$\mathcal{H}_m(\{\mathbf{w}_j\}) = \{f_m(\mathbf{x}, \mathbf{a}) = \frac{1}{m} \sum_{j=1}^m a_j \phi(\mathbf{x}; \mathbf{w}_j)\}.$$

Corresponding function space: reproducing kernel Hilbert space (RKHS) with kernel:

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w} \sim \pi} [\phi(\mathbf{x}; \mathbf{w}) \phi(\mathbf{x}'; \mathbf{w})]$$

Consider the regularized model:

$$\mathcal{L}_{n,\lambda}(\mathbf{a}) = \hat{\mathcal{R}}_n(\mathbf{a}) + \frac{\lambda \|\mathbf{a}\|}{\sqrt{n} \sqrt{m}}, \quad \hat{\mathbf{a}}_n = \operatorname{argmin} \mathcal{L}_{n,\lambda}(\mathbf{a})$$

Up to log terms

$$\mathcal{R}(\hat{\mathbf{a}}_n) \lesssim \frac{\|f^*\|_{\mathcal{H}}^2}{m} + \lambda \frac{\|f^*\|_{\mathcal{H}}}{\sqrt{n}}$$

2. ResNets

$$\begin{aligned}z_{0,L}(\mathbf{x}) &= \mathbf{x}, \\z_{l+1,L}(\mathbf{x}) &= z_{l,L}(\mathbf{x}) + \frac{1}{L} \mathbf{U}_l \sigma \circ (\mathbf{W}_l z_{l,L}(\mathbf{x})), \quad l = 0, 1, \dots, L-1 \\f(\mathbf{x}, \theta) &= \alpha \cdot z_{L,L}(\mathbf{x})\end{aligned}$$

Corresponding function space: “flow-induced function space” (E, Ma and Wu (2019))

Regularized loss function:

$$\mathcal{L}_{n,\lambda}(\theta) = \hat{\mathcal{R}}_n(\theta) + \lambda \|\theta\|_{\mathcal{P}} \sqrt{\frac{2 \log(2d)}{n}}.$$

Up to logarithmic terms, we have:

$$\mathcal{R}(\hat{\theta}) \lesssim \frac{\|f^*\|_{\mathcal{D}}^2}{L} + \lambda \frac{\|f^*\|_{\mathcal{D}}^2}{\sqrt{n}}$$

Variance reduction:

$$\|f^*\|_{\mathcal{D}} \leq \|f^*\|_{\mathcal{B}} \leq \|f^*\|_{\mathcal{H}}$$

3. Multilayer networks

$$f(x) = \sum_{i_L=1}^{m_L} a_{i_L}^L \sigma \left(\sum_{i_{L-1}=1}^{m_{L-1}} a_{i_L i_{L-1}}^{L-1} \sigma \left(\dots \sigma \left(\sum_{i_1=1}^{m_1} a_{i_2 i_1}^1 \sigma \left(\sum_{i_0=1}^{d+1} a_{i_1 i_0}^0 x_{i_0} \right) \right) \right) \right)$$

Corresponding function space: “multilayer space” (E and Wojtowysch (in preparation))

- Rademacher complexity/generalization gap: Same as for Barron space. Let $\mathcal{F}_Q = \{f \in C^0, \|f\|_{\text{multi-layer}} \leq Q\}$. Then $\text{Rad}_S(\mathcal{F}_Q) \leq 2Q \sqrt{\frac{2 \ln(2d+2)}{n}}$
- Inverse approximation theorem holds.
- Direct approximation theorem: holds, but not with Monte-Carlo rate. For f^* in multi-layer space and $m \in \mathbb{N}$, there exists a network f with layers of width $m_\ell = m^{L-\ell+1}$ such that

$$\|f - f^*\|_{L^2(\mathbb{P})} \lesssim \frac{2^L \|f^*\|_{\text{multi-layer}}}{m^{1/(4L-2)}}.$$

Representation of functions for different ML model:

- Random feature model:

$$f(\mathbf{x}) = \mathbb{E}_{\pi}[a(\mathbf{w})\sigma(\mathbf{w}^T \mathbf{x})]$$

- Two-layer neural networks: Barron space

$$f(\mathbf{x}) = \mathbb{E}_{(a,\mathbf{w})\sim\rho}[a\sigma(\mathbf{w}^T \mathbf{x})]$$

- ResNets: Flow-induced space

$$f(\mathbf{x}) = \alpha \cdot \mathbf{z}(1), \dot{\mathbf{z}} = \mathbb{E}_{\{(a,\mathbf{w})\sim\rho_{\tau}\}} \mathbf{a}\sigma(\mathbf{w}^T \mathbf{z}), \mathbf{z}(0) = \mathbf{x}$$

- Multi-layer neural networks: Multi-layer spaces

$$f(\mathbf{x}) = \mathbb{E}_{\theta_L \sim \pi_L} a_{\theta_L}^{(L)} \sigma(\mathbb{E}_{\theta_{L-1} \sim \pi_{L-1}} \dots \sigma(\mathbb{E}_{\theta_1 \sim \pi_1} a_{\theta_2, \theta_1}^1 \sigma(a_{\theta_1}^0 \cdot \mathbf{x})) \dots)$$

What's not known?

- Multi-layer spaces: The approximation error is not optimal. Can it be improved?
- Function spaces for CNNs
- Function space for DenseNets

Optimization:

Does the training process converge to a good solution? How fast?

Generalization:

In particular, is there such thing as “implicit regularization”?

Mean-field formulation

Chizat and Bach (2018), Mei, Montanari and Nguyen (2018), Rotskoff and Vanden-Eijnden (2018), Sirignano and Spiliopoulos (2018)

$$\mathcal{H}_m = \left\{ f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x}) \right\}$$

Let

$$I(\mathbf{u}_1, \dots, \mathbf{u}_m) = \hat{\mathcal{R}}_n(f_m), \quad \mathbf{u}_j = (a_j, \mathbf{w}_j)$$

GD dynamics:

$$\frac{d\mathbf{u}_j}{dt} = -\nabla_{\mathbf{u}_j} I(\mathbf{u}_1, \dots, \mathbf{u}_m), \quad \mathbf{u}_j(0) = \mathbf{u}_j^0, \quad j \in [m]$$

Lemma:. Let

$$\rho(\mathbf{u}, t) = \frac{1}{m} \sum_j \delta_{\mathbf{u}_j(t)}$$

then the GD dynamics described above is equivalent to:

$$\partial_t \rho = \nabla(\rho \nabla V), \quad V = \frac{\delta \hat{\mathcal{R}}_n}{\delta \rho}$$

This is the gradient flow of $\hat{\mathcal{R}}_n$ under the Wasserstein metric.

Convergence to global minima

The functional which is not *displacement convex*, but:

Theorem (Chizat-Bach '18, '20, Wojtowytsch '20)

Let ρ_t be a solution of the Wasserstein gradient flow such that

- ρ_0 has a density on the cone $\Theta := \{|a|^2 \leq |w|^2\}$.
- ρ_0 is omni-directional: Every open cone in Θ has positive measure with respect to ρ_0

Then the following are equivalent.

- 1 The velocity potentials $V = \frac{\delta \mathcal{R}}{\delta \rho}(\rho_t, \cdot)$ converge to a unique limit as $t \rightarrow \infty$.
- 2 $\mathcal{R}(\rho_t) \rightarrow 0$, as $t \rightarrow \infty$.

- 1 There are further technical conditions for the theorem to hold.
- 2 Convergence of subsequences of $\frac{\delta \mathcal{R}}{\delta \rho}(\rho_t, \cdot)$ is guaranteed by compactness.

Training two-layer neural networks under conventional scaling

$$f_m(\mathbf{x}; \mathbf{a}, \mathbf{B}) = \sum_{j=1}^m a_j \sigma(\mathbf{b}_j^T \mathbf{x}) = \mathbf{a}^T \sigma(\mathbf{B}\mathbf{x}),$$

“Xavier-like” initialization:

$$a_j(0) \sim \mathcal{N}(0, \beta^2), \quad \mathbf{b}_j(0) \sim \mathcal{N}(0, I/d)$$

$$\beta = 0 \text{ or } 1/\sqrt{m}$$

The associated *random feature model*: $\{\mathbf{b}_j\}$ frozen, only train $\{a_i\}$

Define Gram matrix $K = (K_{ij})$:

$$K_{i,j} = \frac{1}{n} \mathbb{E}_{\mathbf{b} \sim \pi_0} [\sigma(\mathbf{x}_i^T \mathbf{b}) \sigma(\mathbf{x}_j^T \mathbf{b})].$$

Highly over-parametrized regime

- Good news: Exponential convergence (Du et al (2018))
 - Bad news: Converged solution is no better than that of the random feature model (E, Ma, Wu (2019), Arora et al (2019),)
- Heuristics given in Jacot, Gabriel and Hongler (2018): “neural tangent kernel”

Theorem

Let $\lambda_n = \lambda_{\min}(K)$ and assume $\beta = 0$. Denote by $f_m(\mathbf{x}; \tilde{\mathbf{a}}(t), \mathbf{B}_0)$ the solutions of GD dynamics for the random feature model. For any $\delta \in (0, 1)$, assume that $m \gtrsim n^2 \lambda_n^{-4} \delta^{-1} \ln(n^2 \delta^{-1})$. Then with probability at least $1 - 6\delta$ we have

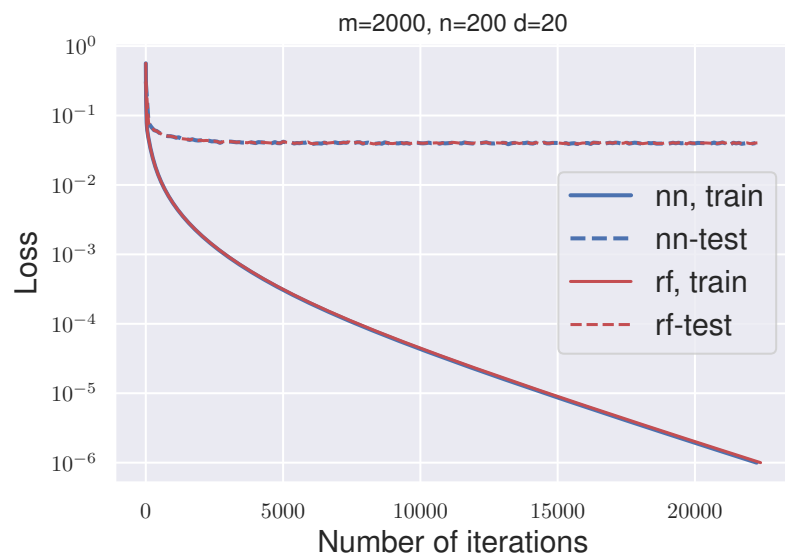
$$\hat{\mathcal{R}}_n(\mathbf{a}(t), \mathbf{B}(t)) \leq e^{-m\lambda_n t} \hat{\mathcal{R}}_n(\mathbf{a}(0), \mathbf{B}(0))$$
$$\sup_{\mathbf{x} \in \mathcal{S}^{d-1}} |f_m(\mathbf{x}; \mathbf{a}(t), \mathbf{B}(t)) - f_m(\mathbf{x}; \tilde{\mathbf{a}}(t), \mathbf{B}_0)| \lesssim \frac{(1 + \sqrt{\ln(\delta^{-1})})^2 \lambda_n^{-1}}{\sqrt{m}}.$$

Disappointing! No “implicit regularization” in this regime, since

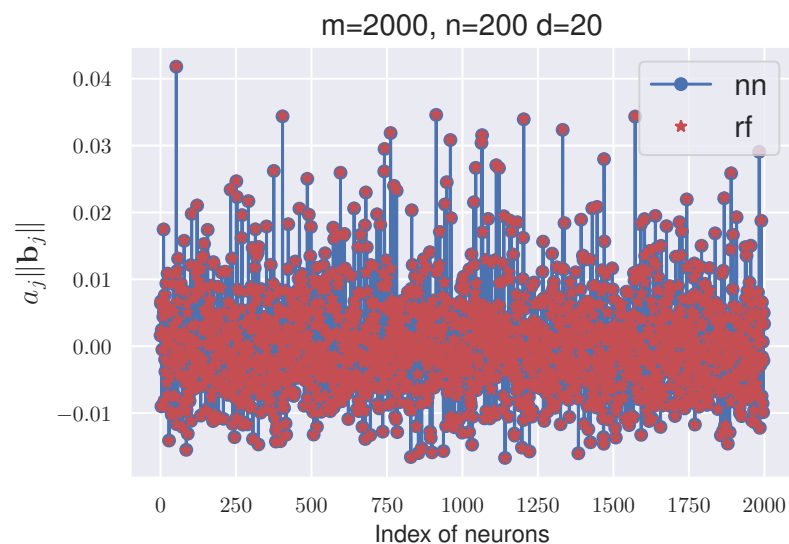
- with explicit regularization, the generalization error is small for all Barron functions.
- without explicit regularization, the generalization error is small only for RKHS functions.

What happens in practice?

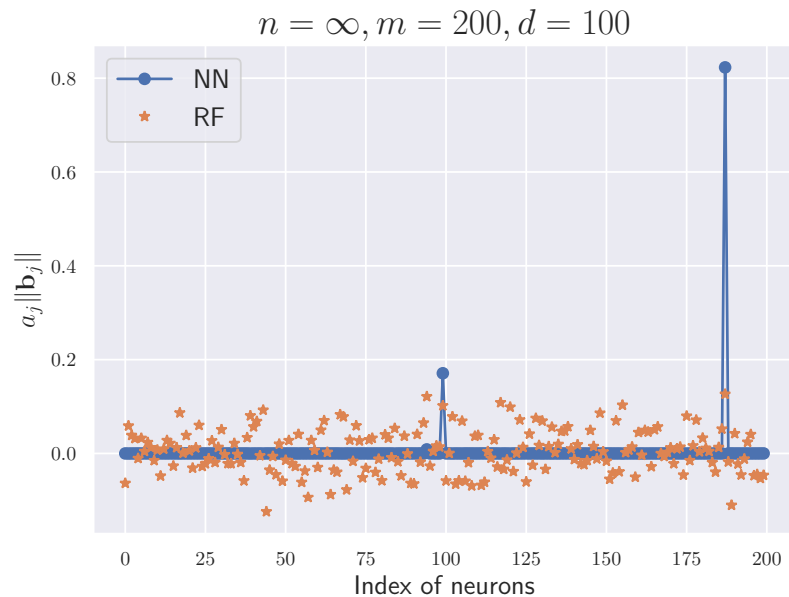
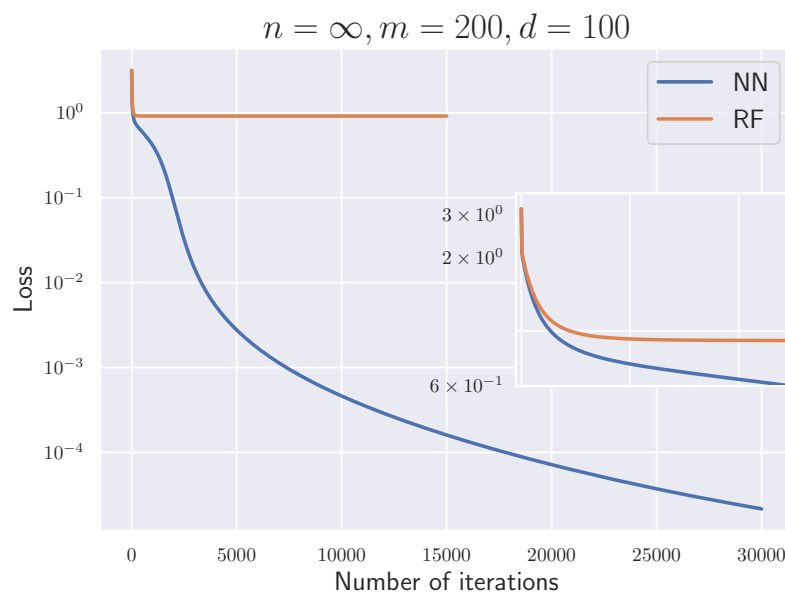
$f^*(\mathbf{x}) = \sigma(\mathbf{e}_1 \cdot \mathbf{x})$ (see Ma, Wu and E (2020) for more results).



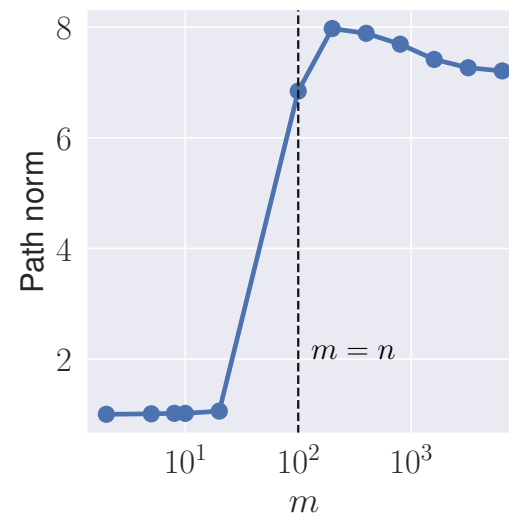
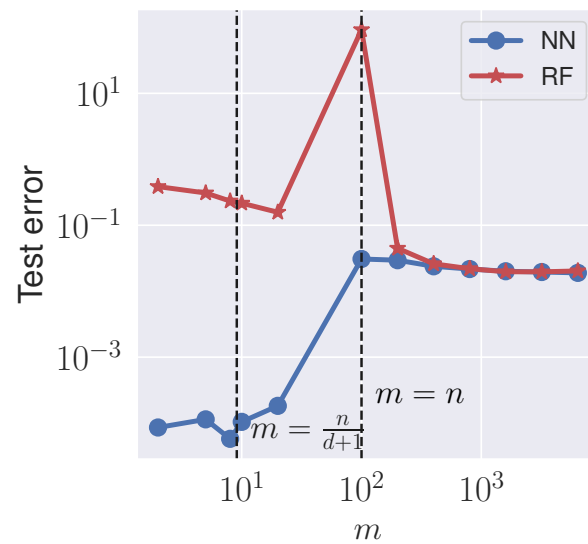
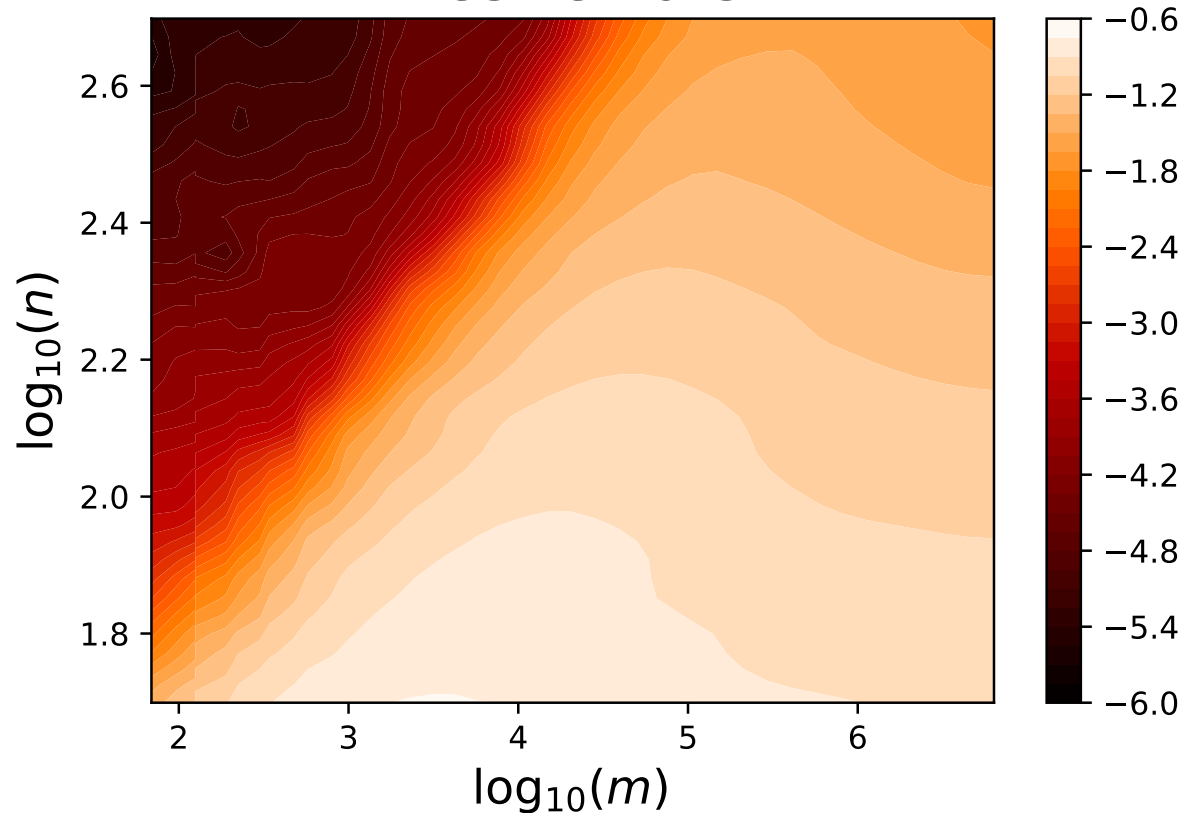
(a)



(b)



Test errors



Neural network-like vs. random feature-like behavior

- Neural network-like behavior
 - two phases: initial random feature-like phase, followed by a second phase with quenching and activation
 - Neurons are divided into two groups: activated ones and background neurons
 - testing error continues to decay after the first phase
- Random feature-like behavior
 - simpler dynamics (e.g. no division into two groups)
 - testing error saturates quickly while training error continues to decay

This is observed for target functions which can be accurately approximated by a small number of neurons (effectively “over-parametrized”).

What we don't know?

Why does GD/SGD converge to good minima?

- mean-field for continuum of neurons: This is a clean analysis problem.
- mean-field for finite neurons: No rigorous results yet.
- conventional scaling: Does there exist a regime with implicit regularization?

“Well-posed” formulations of ML?

Start with a “nice” continuous problem and discretize to get practical ML algorithms

- representation of functions
- the variational problem for minimizing the population risk (loss function)
- gradient flow for the variational problem (training, a PDE-like problem)

Key: The variational problem should be “nice”.

Representation of functions: An illustrative example

Traditional approach:

$$f(\mathbf{x}) = \int_{\mathbb{R}^d} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})} d\boldsymbol{\omega}, \quad f_m(\mathbf{x}) = \frac{1}{m} \sum_j a(\boldsymbol{\omega}_j) e^{i(\boldsymbol{\omega}_j, \mathbf{x})}$$

$\{\boldsymbol{\omega}_j\}$ is a fixed grid, e.g. uniform.

$$\|f - f_m\|_{L^2(X)} \leq C_0 m^{-\alpha/d} \|f\|_{H^\alpha(X)}$$

“New” approach ($\pi =$ probability measure on \mathbb{R}^d):

$$f(\mathbf{x}) = \int_{\mathbb{R}^d} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})} \pi(d\boldsymbol{\omega}) = \mathbb{E}_{\boldsymbol{\omega} \sim \pi} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})}$$

Let $\{\boldsymbol{\omega}_j\}$ be an i.i.d. sample of π .

$$\mathbb{E} \left| f(\mathbf{x}) - \frac{1}{m} \sum_{j=1}^m a(\boldsymbol{\omega}_j) e^{i(\boldsymbol{\omega}_j, \mathbf{x})} \right|^2 = \frac{\text{var}(f)}{m}$$

$\frac{1}{m} \sum_{j=1}^m a(\boldsymbol{\omega}_j) e^{i(\boldsymbol{\omega}_j, \mathbf{x})} =$ two-layer neural network with activation function $\sigma(z) = e^{iz}$.

Integral transform-based representation and the variational problem

Consider the (parametric) representation:

$$f(\mathbf{x}, \theta) = \int_{\mathbb{R}^d} a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x}) \pi(d\mathbf{w}) = \mathbb{E}_{\mathbf{w} \sim \pi} a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x}) = \mathbb{E}_{(a, \mathbf{w}) \sim \rho} a \sigma(\mathbf{w}^T \mathbf{x})$$

θ = parameter:

- $\theta = \{a(\cdot)\}$, π is given
- $\theta = \{a(\cdot), \pi(\cdot)\}$, or equivalently $\theta = \{\rho(\cdot)\}$.

Given a target function f^* , the variational problem for minimizing the population risk:

$$\min_{\theta} \mathcal{R}, \quad \mathcal{R}(\theta) = \mathbb{E}_{\mathbf{x} \sim \mu} (f(\mathbf{x}, \theta) - f^*(\mathbf{x}))^2$$

Conjecture: This is a “nice (convex-like)” variational problem.

Gradient flow for the population risk

Population risk: $\mathcal{R}(f) = \mathbb{E}_{\mathbf{x} \sim \mu} (f(\mathbf{x}) - f^*(\mathbf{x}))^2 = \text{“free energy”}$

$$f(\mathbf{x}) = \int a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x}) \pi(d\mathbf{w}) = \mathbb{E}_{\mathbf{w} \sim \pi} a(\mathbf{w}) \sigma(\mathbf{w}^T \mathbf{x})$$

Follow Halperin and Hohenberg (1977)

- $a = \text{non-conserved}$, use “model A” dynamics (Allen-Cahn):

$$\frac{\partial a}{\partial t} = -\frac{\delta \mathcal{R}}{\delta a}$$

- $\pi = \text{conserved}$ (probability density), use “model B” (Cahn-Hilliard):

$$\frac{\partial \pi}{\partial t} + \nabla \cdot \mathbf{J} = 0$$

$$\mathbf{J} = \pi \mathbf{v}, \quad \mathbf{v} = -\nabla V, \quad V = \frac{\delta \mathcal{R}}{\delta \pi}.$$

Examples

- $\theta = \{a\}$ (fix π , non-conservative).

$$\partial_t a(\mathbf{w}, t) = -\frac{\delta \mathcal{R}}{\delta a}(\mathbf{w}, t) = -\int a(\tilde{\mathbf{w}}, t) K(\mathbf{w}, \tilde{\mathbf{w}}) \pi(d\tilde{\mathbf{w}}) + \tilde{f}(\mathbf{w})$$

$$K(\mathbf{w}, \tilde{\mathbf{w}}) = \mathbb{E}_{\mathbf{x}}[\sigma(\mathbf{w}^T \mathbf{x}) \sigma(\tilde{\mathbf{w}}^T \mathbf{x})], \quad \tilde{f}(\mathbf{w}) = \mathbb{E}_{\mathbf{x}}[f^*(\mathbf{x}) \sigma(\mathbf{w}^T \mathbf{x})]$$

This is an integral equation with a symmetric positive definite kernel.

- $\theta = \{\rho\}$ (conservative)

$$\partial_t \rho = \nabla(\rho \nabla V)$$

$$V(\mathbf{u}) = \frac{\delta \mathcal{R}}{\delta \rho}(\mathbf{u}) = \int \tilde{K}(\mathbf{u}, \tilde{\mathbf{u}}) \rho(d\tilde{\mathbf{u}}) - \tilde{f}(\mathbf{u})$$

This is the same as the mean-field equation.

Discretizing the gradient flows

- Discretizing the population risk (into the empirical risk) using data
- Discretizing the gradient flow
 - particle method – the dynamic version of Monte Carlo
 - smoothed particle method – analog of vortex blob method
 - spectral method – very effective in low dimensions

Particle method for the feature-based model

Continuous problem:

$$\partial_t a(\mathbf{w}, t) = -\frac{\delta \mathcal{R}}{\delta a}(\mathbf{w}) = -\int a(\tilde{\mathbf{w}}, t) K(\mathbf{w}, \tilde{\mathbf{w}}) \pi(d\tilde{\mathbf{w}}) + \tilde{f}(\mathbf{w})$$

$$\pi(d\mathbf{w}) \sim \frac{1}{m} \sum_j \delta_{\mathbf{w}_j}, a(\mathbf{w}_j, t) \sim a_j(t)$$

Discretized version:

$$\frac{d}{dt} a_j(t) = -\frac{1}{m} \sum_k K(\mathbf{w}_j, \mathbf{w}_k) a_k(t) + \tilde{f}(\mathbf{w}_j)$$

This is exactly the GD for the random feature model.

$$f(\mathbf{x}) \sim f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})$$

Discretization of the conservative flow

Consider the integral differential equation (IDE):

$$\partial_t \rho = \nabla(\rho \nabla V)$$

Particle method discretization:

$$\rho(a, \mathbf{w}, t) \sim \frac{1}{m} \sum_j \delta_{(a_j(t), \mathbf{w}_j(t))} = \frac{1}{m} \sum_j \delta_{\mathbf{u}_j(t)}$$

gives rise to

$$\frac{d\mathbf{u}_j}{dt} = -\nabla_{\mathbf{u}_j} I(\mathbf{u}_1, \dots, \mathbf{u}_m)$$

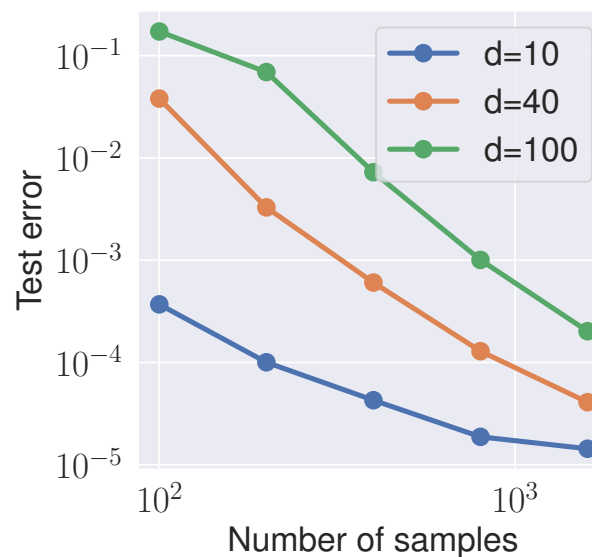
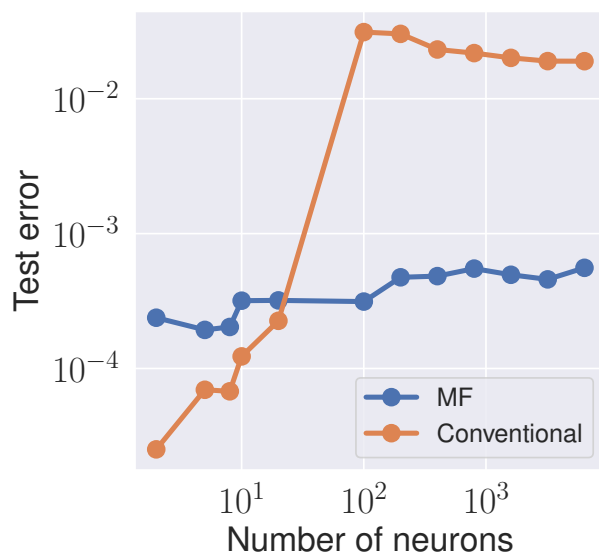
where

$$I(\mathbf{u}_1, \dots, \mathbf{u}_m) = \mathcal{R}(f_m), \quad \mathbf{u}_j = (a_j, \mathbf{w}_j), \quad f_m(\mathbf{x}) = \frac{1}{m} \sum_j a_j \sigma(\mathbf{w}_j^T \mathbf{x})$$

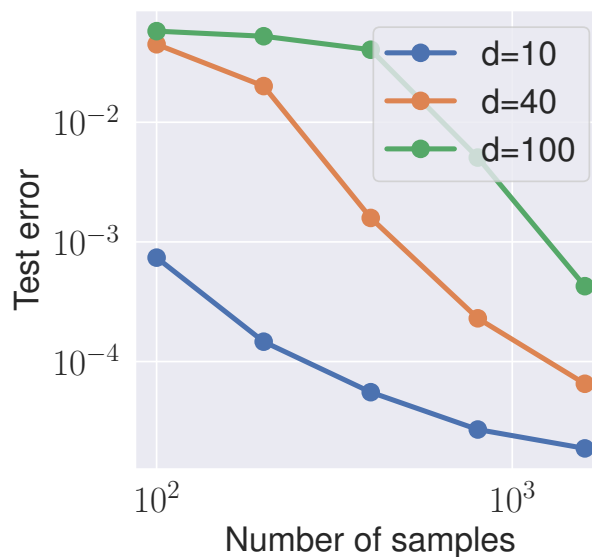
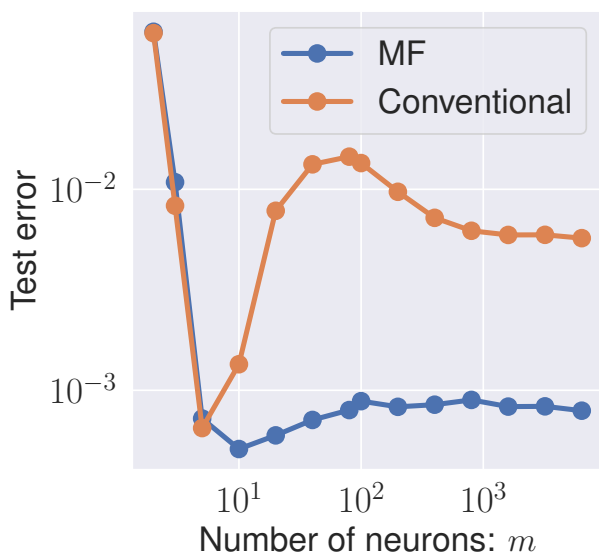
Note that this is exactly the GD dynamics for two-layer neural networks.

Why is “mean-field” or “continuous formulation” better?

Their performance is more robust.



(a) Single neuron f_1^* .



(b) Circle neuron f_2^* .

However, “mean-field” and “continuous” are different viewpoints:

- mean field: discrete \rightarrow continuous by taking the limit (more like interacting particles in stat phys)
- continuous formulation: continuous \rightarrow discrete by discretization (more like the usual numerical analysis situation)

The crucial technical work is in the statics (particularly the representation of functions).

“Continuous” formulation tries to formulate the “first principles” of ML.

It allows us to think “outside the box” about ML.

- Indeed, there are other natural discretizations that lead to “non-neural network-like models” (see E, Ma and Wu (2019)).

ODE viewpoint (E (2017), Haber/Ruthotto (2017), “Neural ODEs” (Chen et al (2018)))

$$\frac{dz}{d\tau} = \mathbf{g}(\tau, \mathbf{z}), \mathbf{z}(0) = \mathbf{x}$$

The flow-map at time 1:

$$\mathbf{x} \rightarrow \mathbf{z}(1)$$

Hypothesis space (or trial functions):

$$f = \alpha^T \mathbf{z}(1)$$

What form of \mathbf{g} should we choose?

A natural choice of g (E, Ma and Wu, 2019):

$$\mathbf{g}(\tau, \mathbf{z}) = \mathbb{E}_{\mathbf{w} \sim \pi_\tau} \mathbf{a}(\mathbf{w}, \tau) \sigma(\mathbf{w}^T \mathbf{z}) = \mathbb{E}_{(\mathbf{a}, \mathbf{w}) \sim \rho_\tau} \mathbf{a} \sigma(\mathbf{w}^T \mathbf{z})$$

where $\{\pi_\tau\}$ or $\{\rho_\tau\}$ is a family of probability distributions.

$$\frac{d\mathbf{z}}{d\tau} = \mathbb{E}_{\mathbf{w} \sim \pi_\tau} \mathbf{a}(\mathbf{w}, \tau) \sigma(\mathbf{w}^T \mathbf{z}) = \mathbb{E}_{(\mathbf{a}, \mathbf{w}) \sim \rho_\tau} \mathbf{a} \sigma(\mathbf{w}^T \mathbf{z})$$

$$f(\mathbf{x}) = \alpha^T \mathbf{z}(\mathbf{x}, 1)$$

Discretize: We obtain the residual neural network model:

$$\mathbf{z}_{l+1} = \mathbf{z}_l + \frac{1}{LM} \sum_{j=1}^M \mathbf{a}_{j,l} \sigma(\mathbf{z}_l^T \mathbf{w}_{j,l}), l = 1, 2, \dots, L-1, \quad \mathbf{z}_0 = V \tilde{\mathbf{x}}$$
$$f_L(\mathbf{x}) = \alpha^T \mathbf{z}_L$$

Particle discretization of the continuous GD flow recovers GD for ResNets.

What we don't know?

Other general representations of functions?

- It has to be some kind of expectation.
- Each such representation may lead to a new kind of ML model.

In a way, neural network models are also very natural: They arise naturally when considering continuous formulations in high dimensions.

Not covered:

- Other approximation results
- Estimates of the generalization gap
- Global minima selection for different optimization algorithms
- Adaptive optimization algorithms (Adam, RMSprop, ...)
- Landscapes, structure of the set of critical points
-

Summary: What do we know and what we don't know?

- approximation/generalization properties of hypothesis space
 - function spaces (RKHS, Barron, flow-induced, multi-layer)– the key is representation of functions
 - generalization error estimates of regularized model
 - CNNs, DenseNets, improvements?
- highly over-parametrized NNs
 - optimizes well
 - generalizes badly
- other regimes under the conventional scaling
 - qualitative behavior (better characterization of NN-like behavior)?
 - behavior for different number of layers?
 - implicit regularization?
- mean-field training dynamics
 - Chizat and Bach and extensions
 - stronger results?
- continuous formulation
 - other representations of functions?
 - other discretizations?
- classification? GAN? RNN? RL?